

Combining Classifiers: Weighted Majority, Bagging, and Stacking

Thursday, 07 November 2002

William H. Hsu

Department of Computing and Information Sciences, KSU

http://www.kddresearch.org

http://www.cis.ksu.edu/~bhsu

Readings: Section 7.5, Mitchell "Bagging, Boosting, and *C4.5*", Quinlan Section 5, "*MLC++* Utilities 2.0", Kohavi and Sommerfield



CIS 732: Machine Learning and Pattern Recognition

Lecture Outline

- Readings
 - Section 7.5, Mitchell
 - Section 5, MLC++ manual, Kohavi and Sommerfield
- This Week's Paper Review: "Bagging, Boosting, and *C4.5*", J. R. Quinlan
- Combining Classifiers
 - Problem definition and motivation: improving accuracy in concept learning
 - General framework: collection of weak classifiers to be improved
- <u>Weighted Majority (WM)</u>
 - Weighting system for collection of algorithms
 - "Trusting" each algorithm in proportion to its training set accuracy
 - Mistake bound for WM
- Bootstrap <u>Aggregating</u> (Bagging)
 - Voting system for collection of algorithms (trained on subsamples)
 - When to expect bagging to work (<u>unstable</u> learners)
- Next Lecture: <u>Boosting</u> the Margin, <u>Hierarchical Mixtures of Experts</u>



CIS 732: Machine Learning and Pattern Recognition

Combining Classifiers

- Problem Definition
 - <u>Given</u>
 - Training data set *D* for supervised learning
 - D drawn from common instance space X
 - Collection of inductive learning algorithms, hypothesis languages (inducers)
 - Hypotheses produced by applying inducers to *s*(*D*)
 - s: X vector \rightarrow X' vector (sampling, transformation, <u>partitioning</u>, etc.)
 - Can think of hypotheses as definitions of <u>prediction algorithms</u> ("classifiers")
 - <u>Return</u>: new prediction algorithm (*not* necessarily \in *H*) for $x \in X$ that combines outputs from collection of prediction algorithms
- Desired Properties
 - Guarantees of performance of combined prediction
 - e.g., mistake bounds; ability to improve weak classifiers
- Two Solution Approaches
 - Train and apply each inducer; learn combiner function(s) from result
 - Train inducers and combiner function(s) concurrently



Principle: Improving Weak Classifiers



CIS 732: Machine Learning and Pattern Recognition

Framework: Data Fusion and Mixtures of Experts

- What *Is* A <u>Weak</u> Classifier?
 - One not guaranteed to do better than random guessing (1 / number of classes)
 - Goal: combine multiple weak classifiers, get one at least as accurate as strongest
- Data Fusion
 - Intuitive idea
 - Multiple sources of data (sensors, domain experts, etc.)
 - Need to combine systematically, plausibly
 - Solution approaches
 - Control of intelligent agents: <u>Kalman filtering</u>
 - General: <u>mixture estimation</u> (sources of data \Rightarrow predictions to be combined)
- Mixtures of Experts
 - Intuitive idea: "experts" express hypotheses (drawn from a hypothesis space)
 - Solution approach (next time)
 - <u>Mixture model</u>: estimate <u>mixing coefficients</u>
 - Hierarchical mixture models: divide-and-conquer estimation method



Weighted Majority: Idea

- Weight-Based Combiner
 - <u>Weighted votes</u>: each prediction algorithm (classifier) h_i maps from $x \in X$ to $h_i(x)$
 - Resulting prediction in set of legal class labels
 - NB: as for <u>Bayes</u> Optimal <u>Classifier</u>, resulting *predictor* not necessarily in H
- Intuitive Idea
 - Collect votes from pool of prediction algorithms for each training example
 - Decrease weight associated with each algorithm that guessed wrong (by a multiplicative factor)
 - Combiner predicts <u>weighted majority</u> label
- Performance Goals
 - Improving training set accuracy
 - Want to combine weak classifiers
 - Want to bound number of mistakes in terms of minimum made by any one algorithm
 - Hope that this results in good generalization quality



CIS 732: Machine Learning and Pattern Recognition

Weighted Majority: Procedure

- Algorithm *Combiner-Weighted-Majority* (*D*, *L*)
 - $n \leftarrow L.size$
 - $m \leftarrow D.size$
 - FOR *i* ← 1 TO *n* DO
 - $P[i] \leftarrow L[i]$. Train-Inducer (D)
 - *w_i* ← 1
 - FOR $j \leftarrow 1$ TO m DO
 - $q_0 \leftarrow 0, q_1 \leftarrow 0$
 - FOR *i* ← 1 TO *n* DO

IF P[i](D[j]) = 0 THEN $q_0 \leftarrow q_0 + w_i$

IF P[i](D[j]) = 1 THEN $q_1 \leftarrow q_1 + w_i$

Prediction[*i*][*j*] ← $(q_0 > q_1)$? 0 : $((q_0 = q_1)$? *Random* (0, 1): 1)

IF *Prediction[i][j] ≠ D[j].target* THEN

 $w_i \leftarrow \beta w_i$

- RETURN Make-Predictor (w, P)

```
// number of inducers in pool
// number of examples \langle x \equiv D[i], c(x) \rangle
```

// P[i]: ith prediction algorithm
// initial weight
// compute WM label

// vote for 0 (-) // else vote for 1 (+)

// *c*(*x*) ≡ *D*[*j*].*target*

// β < 1 (i.e., penalize)



Weighted Majority: Properties

- Advantages of WM Algorithm
 - Can be adjusted incrementally (without retraining)
 - Mistake bound for WM
 - Let *D* be any sequence of training examples, *L* any set of inducers
 - Let k be the minimum number of mistakes made on D by any L[i], $1 \le i \le n$
 - <u>Property</u>: number of mistakes made on *D* by *Combiner-Weighted-Majority* is at most 2.4 (*k* + lg *n*)
- Applying *Combiner-Weighted-Majority* to Produce Test Set Predictor
 - *Make-Predictor*: applies abstraction; returns <u>funarg</u> that takes input $x \in D_{test}$
 - Can use this for incremental learning (if c(x) is available for new x)
- Generalizing Combiner-Weighted-Majority
 - Different input to inducers
 - Can add an argument s to sample, transform, or partition D
 - Replace $P[i] \leftarrow L[i]$. Train-Inducer (D) with $P[i] \leftarrow L[i]$. Train-Inducer (s(i, D))
 - Still compute weights based on performance on D
 - Can have q_c ranging over more than 2 class labels



Bagging: Idea

- <u>Bootstrap Agg</u>regating aka Bagging
 - Application of <u>bootstrap sampling</u>
 - <u>Given</u>: set *D* containing *m* training examples
 - Create S[i] by drawing m examples at random with replacement from D
 - S[i] of size m: expected to leave out 0.37 of examples from D
 - Bagging
 - Create *k* bootstrap samples *S*[1], *S*[2], ..., *S*[*k*]
 - Train distinct inducer on each S[i] to produce k classifiers
 - Classify new instance by classifier vote (equal weights)
- Intuitive Idea
 - "Two heads are better than one"
 - Produce multiple classifiers from one data set
 - NB: same inducer (multiple instantiations) or different inducers may be used
 - Differences in samples will "smooth out" sensitivity of L, H to D



Bagging: Procedure

- Algorithm Combiner-Bootstrap-Aggregation (D, L, k)
 - FOR $i \leftarrow 1$ TO k DO
 - S[i] ← Sample-With-Replacement (D, m)
 - Train-Set[i] \leftarrow S[i]
 - *P*[*i*] ← *L*[*i*].*Train-Inducer* (*Train-Set*[*i*])
 - RETURN (Make-Predictor (P, k))
- Function Make-Predictor (P, k)
 - **RETURN** (fn $x \Rightarrow$ *Predict* (*P*, *k*, *x*))
- Function Predict (P, k, x)
 - FOR *i* ← 1 TO *k* DO
 - $Vote[i] \leftarrow P[i](x)$
 - RETURN (argmax (Vote[i]))
- Function Sample-With-Replacement (D, m)
 - RETURN (*m* data points sampled i.i.d. uniformly from *D*)



Bagging: Properties

- Experiments
 - [Breiman, 1996]: Given sample S of labeled data, do 100 times and report average
 - 1. Divide S randomly into test set D_{test} (10%) and training set D_{train} (90%)
 - 2. Learn decision tree from D_{train}
 - $e_{S} \leftarrow \text{error of tree on } T$
 - 3. Do 50 times: create bootstrap *S*[*i*], learn decision tree, prune using *D*

 $e_B \leftarrow$ error of majority vote using trees to classify T

- [Quinlan, 1996]: Results using UCI <u>Machine Learning Database Repository</u>
- When Should This Help?
 - When learner is <u>unstable</u>
 - Small change to training set causes large change in output hypothesis
 - True for decision trees, neural networks; not true for *k*-nearest neighbor
 - Experimentally, bagging can help substantially for unstable learners, can somewhat degrade results for stable learners



Bagging: Continuous-Valued Data

- Voting System: Discrete-Valued Target Function Assumed
 - Assumption used for WM (version described here) as well
 - Weighted vote
 - Discrete choices
 - Stacking: generalizes to continuous-valued targets *iff* combiner inducer does
- Generalizing Bagging to Continuous-Valued Target Functions
 - Use mean, not mode (aka argmax, majority vote), to combine classifier outputs
 - Mean = expected value
 - $\phi_A(x) = \mathsf{E}_D[\phi(x, D)]$
 - $\phi(x, D)$ is base classifier
 - $\phi_A(x)$ is aggregated classifier
 - $(\mathsf{E}_D[y \phi(x, D)])^2 = y^2 2y \cdot \mathsf{E}_D[\phi(x, D)] + \mathsf{E}_D[\phi^2(x, D)]$
 - Now using $E_D[\phi(x, D)] = \phi_A(x)$ and $EZ^2 \ge (EZ)^2$, $(E_D[y \phi(x, D)])^2 \ge (y \phi_A(x))^2$
 - Therefore, we expect lower error for the bagged predictor ϕ_A



Stacked Generalization: Idea

- Stacked Generalization aka <u>Stacking</u>
- Intuitive Idea



CIS 732: Machine Learning and Pattern Recognition

Stacked Generalization: Procedure

- Algorithm *Combiner-Stacked-Gen* (*D*, *L*, *k*, *n*, *m*', *Levels*)
 - Divide *D* into *k* segments, *S*[1], *S*[2], ..., *S*[*k*]
 - FOR *i* ← 1 TO *k* DO
 - Validation-Set ← S[i]
 - FOR *j* ← 1 TO *n* DO
 Train-Set[*j*] ← Sample-With-Replacement (*D* ~ S[*i*], *m*²) // *m m*/k examples
 - IF *Levels* > 1 THEN
 - $P[j] \leftarrow Combiner-Stacked-Gen (Train-Set[j], L, k, n, m', Levels 1)$

ELSE

 $P[j] \leftarrow L[j]$. Train-Inducer (Train-Set[j])

• Combiner ← L[0]. Train-Inducer (Validation-Set.targets,

Apply-Each (P, Validation-Set.inputs))

- Predictor ← Make-Predictor (Combiner, P)
- RETURN Predictor
- Function Sample-With-Replacement: Same as for Bagging



// Assert D.size = m

// Base case: 1 level

// *m/k* examples

Stacked Generalization: Properties

- Similar to Cross-Validation
 - *k*-fold: rotate validation set
 - Combiner mechanism based on validation set as well as training set
 - Compare: <u>committee-based combiners</u> [Perrone and Cooper, 1993; Bishop, 1995] aka <u>consensus under uncertainty / fuzziness</u>, <u>consensus models</u>
 - Common application with cross-validation: treat as overfitting control method
 - Usually improves generalization performance
- Can Apply Recursively (Hierarchical Combiner)
 - Adapt to inducers on different subsets of input
 - Can apply s(Train-Set[j]) to transform each input data set
 - e.g., attribute partitioning [Hsu, 1998; Hsu, Ray, and Wilkins, 2000]
 - Compare: <u>Hierarchical Mixtures of Experts (HME</u>) [Jordan *et al*, 1991]
 - Many differences (validation-based vs. mixture estimation; online vs. offline)
 - Some similarities (hierarchical combiner)



Other Combiners

- So Far: <u>Single-Pass</u> Combiners
 - First, train each inducer
 - Then, train combiner on their output and evaluate based on criterion
 - Weighted majority: training set accuracy
 - Bagging: training set accuracy
 - Stacking: validation set accuracy
 - **<u>Finally</u>**, apply combiner function to get new prediction algorithm (classfier)
 - Weighted majority: weight coefficients (penalized based on mistakes)
 - Bagging: voting committee of classifiers
 - Stacking: validated hierarchy of classifiers with trained combiner inducer
- Next: <u>Multi-Pass</u> Combiners
 - Train inducers and combiner function(s) concurrently
 - Learn how to *divide* and *balance* learning problem across multiple inducers
 - Framework: mixture estimation

Terminology

- Combining Classifiers
 - <u>Weak classifiers</u>: not guaranteed to do better than random guessing
 - <u>Combiners</u>: functions *f*: *prediction vector* \times *instance* \rightarrow *prediction*
- Single-Pass Combiners
 - <u>Weighted Majority (WM)</u>
 - Weights prediction of each inducer according to its training-set accuracy
 - Mistake bound: maximum number of mistakes before converging to correct h
 - Incrementality: ability to update parameters without complete retraining
 - <u>Bootstrap Aggregating</u> (aka <u>Bagging</u>)
 - Takes vote among multiple inducers trained on different samples of D
 - <u>Subsampling</u>: drawing one sample from another (*D* ~ *D*)
 - Unstable inducer: small change to D causes large change in h
 - Stacked Generalization (aka Stacking)
 - Hierarchical combiner: can apply recursively to re-stack
 - Trains <u>combiner inducer</u> using validation set



Summary Points

- Combining Classifiers
 - Problem definition and motivation: improving accuracy in concept learning
 - General framework: collection of weak classifiers to be improved (data fusion)
- <u>Weighted Majority (WM)</u>
 - Weighting system for collection of algorithms
 - Weights each algorithm in proportion to its training set accuracy
 - Use this weight in performance element (and on test set predictions)
 - Mistake bound for WM
- <u>Bootstrap Aggregating (Bagging)</u>
 - Voting system for collection of algorithms
 - Training set for each member: sampled with replacement
 - Works for unstable inducers
- Stacked Generalization (aka Stacking)
 - Hierarchical system for combining inducers (ANNs or other inducers)
 - Training sets for "leaves": sampled with replacement; combiner: validation set
- Next Lecture: <u>Boosting</u> the Margin, <u>Hierarchical Mixtures of Experts</u>

